

`{acuselton,dushizima}@lbl.gov`

Copyright is held by the author/owner(s).
SC'11 Companion, November 12–18, 2011, Seattle, Washington, USA.
 ACM 978-1-4503-1030-7/11/11.

The long term average data rate to or from a Lustre server (OST) is about 10MB/s. Only about 5% of all observations are above this average, that is, the observations are “bursty.” An “event” is a contiguous sequence of observations that contain points above the long term average. Figure 1 illustrates the event detection concept, where the time series can be segmented into a sequence of events, with each characterized by the tuple (start time, duration, size). The total amount of data moved in such events accounts for the majority of all I/O to or from the file system.

3. STATISTICS OF EVENTS

After calculating histograms of events organized by length, size, and rate, we noticed that the vast majority of events are transient, occupying only a single observation. The average length of an event is 2 observations (ten seconds), but half of all I/O takes place in events lasting more than four observations. For the smallest events there is an exponential decline in prevalence with increased size. The histograms show that there are several peaks in the distribution for larger values of length, size or rate. Each peak represents a preferred mode of I/O, for example size 2GB is much more likely than slightly larger or smaller sizes. The prevalence of events in the peak indicates a characteristic of a particular application’s I/O pattern or of the workload.

We notice that a large fraction of all I/O is taking place well below the optimum rate of 400MB/s. This fact has implications for the design of applications generating the I/O workload. Poorly organized I/O may use up all the available bandwidth while still achieving suboptimal performance - a competing transaction would not find the extra bandwidth available. To the extent that this sub-optimality is unavoidable it needs to be factored into the design of the I/O subsystem. A job running its I/O at 100MB/s represents four times the opportunity cost compared to running at the maximum available I/O rate.

4. PRELIMINARY RESULTS

Our preliminary findings aimed at mapping the signals into features that can be used to group events. Also, we designed graphical representations as those in Figures 2, which present 48 time series for the Franklin /scratch read (respectively write) values as spark-lines. The x-axis represents 24 hours, and the y-axis gives one spark line for each OST. Each segment indicates a transaction’s duration, with its width proportional to the amount of I/O. This visualization clearly reflects the coordinated timing of some activity across the entire I/O subsystem. This is a fast and easy way to identify correlations between simultaneous events across OSTs. The framework uses the R Statistical tools, particularly the package *multicore* for parallel calculation of transactions given multiple OSTs, and package *fields* for appropriate color to be associated with different OSTs.

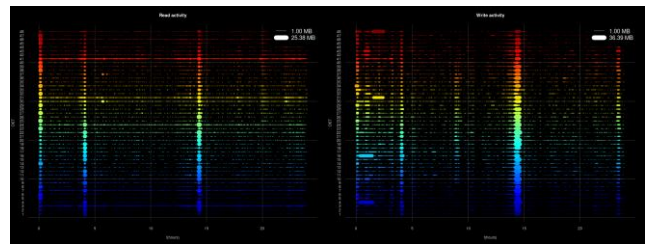


Figure 2. Correlation among OSTs: read and write activity on Franklin using sparklines to indicate I/O transfer volume.

The present analysis identifies events by their timing, duration, and the amount of I/O involved [4]. Those characteristics can often emphasize the I/O of a particular application. This scheme allows calculation of metrics on a continuous basis, allowing dynamic recognition of changes in the workload. The identification still requires a human to search for patterns. Future releases will address algorithms and information visualization regarding the likelihood of a particular user to be associated with transactions which performed large amounts of I/O, e.g. clustering of users based on typical I/O transactions. Another possibility is to identify events that are periodic, as appears to be the case in several data entries.

5. ACKNOWLEDGMENTS

This work was supported by the Director, Office of Science of the U.S. Department of Energy under Contract No. DE-AC02-05CH11231.

6. REFERENCES

- [1] E. Im, I. Bustany, C. Ashcraft, J. Demmel, K. A. Yelick. Performance Tuning of Matrix Triple Products Based on Matrix Structure. In Proceedings of PARA'2004. pp.740-746, 2004.
- [2] C. M. Herb Wartens, Jim Garlick. LMT - The Lustre Monitoring Tool. <http://code.google.com/p/lmt>
- [3] A. Uselton, K. Antypas, D. Ushizima, J. Sukharev, "File System Monitoring as a Window Into User I/O Requirements", CUG-2010, Edinburgh, UK, May 24-27th, 2010.
- [4] D. Ushizima, A. Uselton, K. Antypas, J. Sukharev "Minimizing I/O contention at NERSC using data analysis", Workshop on Algorithms for Modern Massive Data Sets (MMDS'10), Stanford, CA, June 15-18, 2010